

Image Tamper Detection and Recovery by Intersecting Signatures

Chun-Hung Chen¹, Yuan-Liang Tang^{*2}, Wen-Shyong Hsieh³, Min-Shiang Hwang⁴

^{1,3} Department of Computer Science and Engineering, National Sun Yat-sen University, Taiwan

² Department of Information Management, Chaoyang University of Technology, Taiwan

³ Department of Computer and Communication, Shu-Te University, Taiwan

⁴ Department of Computer Science and Information Engineering, Asia University, Taiwan

*Corresponding author, e-mail: yltang@cyut.edu.tw

Abstract

In this paper, we propose an exact image authentication scheme that can, in the best case, detect image tampering with the accuracy of one pixel. This method is based on constructing blocks in the image in such a manner that they intersect with one another in different directions. Such a technique is very useful to identify whether an individual image pixel has been tampered with. Moreover, the tampered region can be well recovered with the embedded recover data.

Keywords: information security, digital watermarking, image authentication, digital signatures

1. Introduction

As digital technologies advance, more and more publications are produced in digital formats and transmitted via the Internet. Accompanying such advance, however, unauthorized use, illegal copying, and malicious modification of digital products have become serious problems. Researchers thus try to find various ways to protect digital products; solutions include copyright assertion, content authentication, etc. In the area of image content authentication, the integrity of an image is regarded very important and must therefore be realized. A common approach is the use of digital watermarking techniques. Digital watermarking serves many purposes, for example, proof of ownership, content authentication, copy control, and so on.

Researchers have developed various image authentication techniques to detect if an image has experienced unauthorized modification. Some of them can only detect whether the image as a whole has been altered. Others may have the additional capability to detect if a certain part of the image has been tampered with. Liu *et al.* [1] studied the Zernike moment values which are generated from low DWT subbands. They found that the quantized values are robust to common processing operations but fragile to malicious attacks. Therefore, they embedded the watermark by quantizing the Zernike moment values, and the locations (i.e., blocks) suffered from malicious attacks can be identified through examining the extracted values. Their method has moderate robustness against JPEG compression. In Rawat and Raman's scheme [2], two chaotic maps are used in order to enhance the security of the watermarked images. The pixels in the image are disturbed using the first chaotic map and are further separated into bit planes with the least significant bit used for watermark embedding. A binary watermark is scrambled by the second chaotic map. The watermarked images can avoid counterfeiting attacks. Xi'an [3] scrambled a bi-level watermark by the Arnold transform, and the Human Visual System is used to determine the quantization step. The scrambled watermark is then inserted into the low DWT coefficients. Tamper areas can then be localized by comparing the extracted and the original watermarks. Patra *et al.* [4] convert the images into the DCT domain and quantize the low-frequency coefficients according to the target levels determined by the Chinese Remainder Theorem. Their method is computationally efficient and is able to withstand such attacks as JPEG compression, sharpening, and brightening. Qi *et al.* [5] used two content-based watermarks to protect the images. One of them is generated by an edge detector for the purpose of detecting tiny changes, and the other is generated from the relationship between the wavelet coefficients for localizing tampered regions. Both watermarks are embedded into middle- and high-frequency DWT coefficients. Finally, the generated watermarks and extracted watermarks are compared to authenticate the image, and a malicious

attack is identified if error pixels are clustered together. Their method is robust against several image processing operations, including JPEG compression. In Wu's work [6], the image is divided into blocks, and all hashes derived from the MSBs of each block are further encoded using an error correcting code (ECC). The parities, rather than the codewords, of the ECC are separated and embedded into the LSBs of each block. During authentication, the original hashes can be recovered if the number of tampered blocks is less than a threshold. The hash of a block is produced and compared with the original to identify if the block is tampered with. This method has fine granularity on detecting tampered regions.

Although many techniques have been proposed for image authentication, most of them can only detect if an image or part of it, as a whole, is modified; very few can identify image tampers down to the granularity of one-pixel level. In some applications, such ability could be extremely essential. For example, if an image is used as a critical piece of evidence in the court or in a police investigation, a generalized answer as to whether the entire image or part of it is altered to some degree may not be acceptable by the law. To be exact, it may be mandatory that the image should not allow for even tiny modification ever since it was taken. Our previous work [7] is able to achieve such exact authentication. In this paper, we further propose an authentication scheme that is able not only to detect and locate image tampers with the accuracy of one pixel at the best case, but also to recover the tampered data. This is done by first constructing linear blocks in the image in such a way that they intersect with one another in different directions. Second, a signature is created for each block for the purpose of authentication. Third, the signature is embedded back into the image in order to protect each pixel by the four signatures and any tampered pixel can be pinpointed by examining its corresponding signatures. And finally, the quantized DCT coefficients are generated and embedded for the purpose of recovering the tamper regions. The rest of the paper is organized as follows. The proposed technique is described in Section 2, followed by experimental results in Section 3. Section 4 presents a security analysis. A comparison of detection granularity is shown in Section 5. Finally, Section 6 gives some concluding remarks.

2. Proposed Technique

2.1. Constructing the Authentication, Signature, and Recovery Blocks

Without loss of generality, we assume that 8-bit grayscale images are dealt with. For other formats, the same technique applies, too. The image is first divided into equal-sized $B \times B$ blocks, which are referred to as the *authentication blocks*. And then, in each authentication block, four sets of pixels are collected in four directions: horizontal, vertical, -45° and 45° wrap-around diagonals. These sets of linear blocks are referred to as *signature blocks*. Namely,

$$\begin{aligned} H_i &= \{p_{ij} \mid j = 0, \dots, B-1\}: \text{horizontal blocks, } i = 0, \dots, B-1, \\ V_j &= \{p_{ij} \mid i = 0, \dots, B-1\}: \text{vertical blocks, } j = 0, \dots, B-1, \\ X_m &= \{p_{ij} \mid i = 0, \dots, B-1, j = (m+i) \bmod B\}: -45^\circ \text{ blocks, } m = 0, \dots, B-1, \text{ and} \\ Y_n &= \{p_{ij} \mid i = (B+n-j) \bmod B, j = 0, \dots, B-1\}: 45^\circ \text{ blocks, } n = 0, \dots, B-1, \end{aligned}$$

where p_{ij} denotes the image pixel and \bmod is the modulo operation. Figure 1 depicts such a construction.

authentication and the recovery blocks are 160×160 (i.e., $B = 160$) and 10×10 (i.e., $R = 10$), respectively.

When performing authentication, the construction of the authentication blocks and signature blocks are repeated as before, followed by the same DES encryption process. Now, the results can be matched against those extracted from the LSBs in the image. The mismatched blocks will be recorded in the following four sets, respectively:

$$\begin{aligned} E_H &= \{H_i \mid i = i_0, i_1, \dots, i_{h-1}\}: \text{horizontal blocks,} \\ E_V &= \{V_j \mid j = j_0, j_1, \dots, j_{v-1}\}: \text{vertical blocks,} \\ E_X &= \{X_m \mid m = m_0, m_1, \dots, m_{x-1}\}: -45^\circ \text{ blocks, and} \\ E_Y &= \{Y_n \mid n = n_0, n_1, \dots, n_{y-1}\}: 45^\circ \text{ blocks,} \end{aligned}$$

where h , v , x , and y are the respective numbers of mismatched blocks. The basic idea of the algorithm is that since each pixel is protected by four signatures and the signature blocks intersect with one another, if a specific pixel is indeed tampered with, mismatches will occur in all of its four corresponding signatures. On the contrary, if some of the corresponding signatures are matched, it can be concluded that the pixel has not been altered. Fig. 3 illustrates the algorithm of tamper detection in an authentication block. Each pixel, p_{ij} ($0 \leq i, j \leq B-1$), in a block is checked to see if it is tampered with. This is done by examining its corresponding four signatures, i.e., horizontal, vertical, and two diagonal ones. If all four signatures mismatch, the pixel will be reported as been tampered with.

```

for ( $i = 0, \dots, B-1$ )
  if ( $H_i \in E_H$ ) # Horizontal signature mismatches
    for ( $j = 0, \dots, B-1$ )
      if ( $V_j \in E_V$ ) # Vertical signature mismatches
         $m \leftarrow j-i$ ;
        if ( $m < 0$ )
           $m \leftarrow m+B$ ; # Wrap around
        end if
         $n \leftarrow i+j$ ;
        if ( $n > B-1$ )
           $n \leftarrow n-B$ ; # Wrap around
        end if
        if ( $m \in E_X$ ) and ( $n \in E_Y$ )
          # Both diagonal signatures mismatch
           $p_{ij}$ : tampered pixel;
        end if
      end if
    end for
  end if
end for

```

Figure 3. The algorithm of tamper detection

If a mismatch is detected, the recovery is performed using a reference image. The construction of the reference image starts by extracting the recovery data from each recovery data. If all of pixels of a recovery block are authentic, the embedded recovery data are considered as *valid*. If the data extracted from the corresponding block are invalid, the recovery block will be constructed with its valid data using de-quantization and inverse DCT, followed by scaling back the block size from 8×8 to $R \times R$. The reference image is obtained after all of the blocks are produced, and the recovery is done by replacing the tampered pixels with the

corresponding pixels of the reference image if the corresponding pixels are constructed by the valid data.

2.2. Analysis

One of the main shortcomings of most other tamper detection techniques is that they usually create a signature for an image block, and if a mismatch occurs, the block as a whole is identified as being tampered with. There is no way of distinguishing which pixel (or pixels) is the victim. The essence of the proposed scheme lies on the fact that each pixel is protected by four intersecting signature blocks. Whenever one pixel is tampered with, it causes the four corresponding signatures to mismatch and, through the intersecting structure, the tampered pixel can be easily pinpointed. In other words, if less than four mismatches occur for a pixel, we can eliminate the possibility of tampering. This method is thus very accurate in identifying tampered pixels as well as their locations in the image. There are, however, some conditions in which this scheme will make false positive reports. Figure 4 illustrates such a situation, in which the black pixels represent those pixels that have been altered by attackers. The four sets of mismatched blocks are: $E_H = \{H_{i1}, H_{i2}, H_{i3}\}$, $E_V = \{V_{j1}, V_{j2}, V_{j3}\}$, $E_X = \{X_{m1}, X_{m2}, X_{m3}\}$, and $E_Y = \{Y_{n1}, Y_{n2}, Y_{n3}\}$. It is obvious that, besides the four black pixels, the system will erroneously report the center pixel (represented by a white pixel) as a tampered one, i.e., a false positive.

Actually, the number of tampered pixels (or the size of the tampered region) determines the number of mismatch signatures. If the former increases, the latter increases, too. There is no constraint on the maximum size of a tampered region; however, the shape of the region does affect the number of false positives in the detection. Figure 5 illustrates an example, in which the black pixels represent the tampered pixels. As our algorithm identifies tampered pixels by four intersecting signatures, the set of reported pixels will form a convex shape (the red polygon in the figure). As all pixels in the convex shape are reported as tampered with, the unchanged pixels (white pixels) are false positives. The situation gets worse if the tampered pixels spread randomly across the image. The result of a simulation is shown Figure 6, in which the tampered pixels are generated randomly across the authentication block (160×160). It can be seen that the number of false positives increases rapidly as the number of tampered pixels increases, almost reaching 16,000 when the latter is only a few hundreds. Such a phenomenon verifies the above analysis that as the reported pixels form a convex shape, if the locations of the tampered pixels are randomly generated, the convex area will become very large, which results in a great number of false positives. Denoting the numbers of reported, tampered, and false positive pixels as R , T , and F , respectively, the following relationship holds:

$$R = T + F$$

Therefore, when R achieves its highest (i.e., the size of the block), increasing T will certainly decrease F , which explains the peak in Figure 6. The same false positive effect is also expected in other existing block-based authentication techniques. In practice, however, as an attacker usually tries to alter the semantics of the image, tampered pixels tend to cluster together (may be in several locations). Randomly altering the pixels is meaningless and hence not likely to happen.

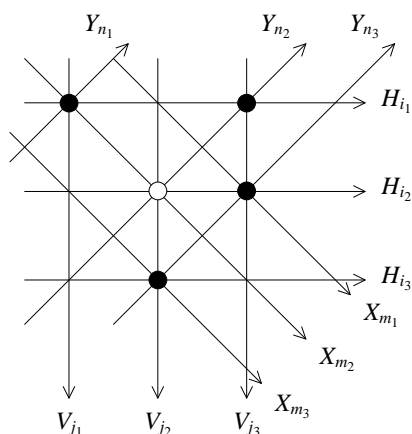


Figure 4. Example of a false positive

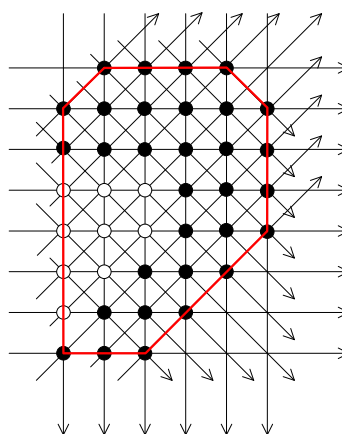


Figure 5. A tampered region (black pixels)

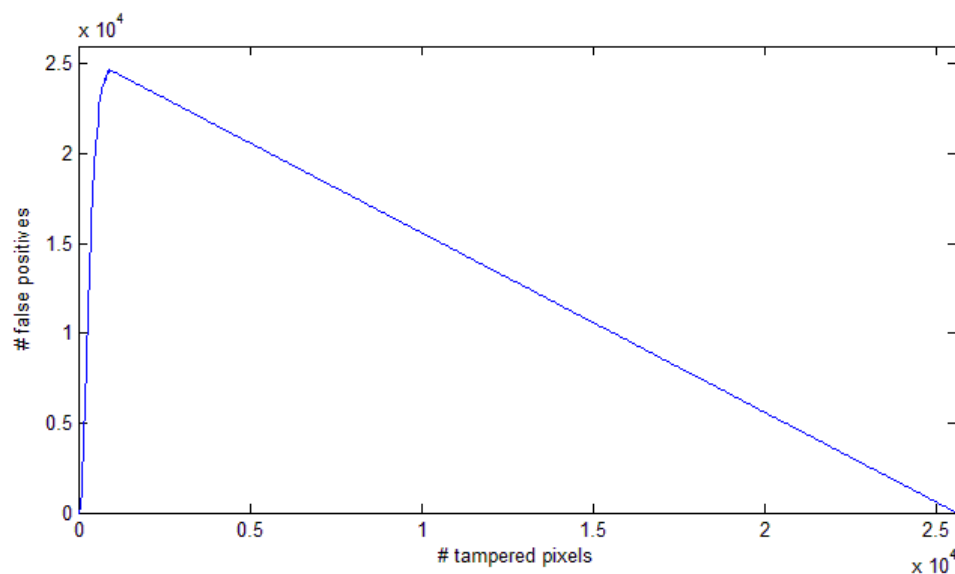


Figure 6. Number of false positives vs. number of randomly tampered pixels

2.3. Handling Irregular Image Sizes

If the image size is not multiples of the size of the authentication block, something must be done for the extra areas. Since 64 bits are required for a signature, every 32 pixels can be collected to form an individual authentication block. In those areas, however, if the signature mismatches, it can be only concluded that one or more pixels in such a block could have been altered. Figure 7 shows such a situation. The granularity of identification in those areas is thus 32 pixels.

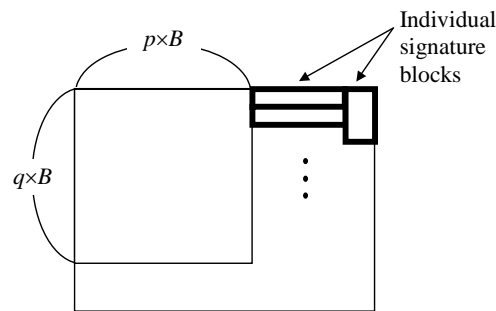


Figure 7. Handling the situation in which the image size is not multiples of $B \times B$

3. Experimental Results

A set of 480×480 grayscale images (*Airplane*, *Baboon*, *Lena*, and *Peppers*) were used to test the proposed algorithm: Figure 8(a) and 8(b) show the original images and the corresponding data embedded versions. The PSNR between the original and the embedded images are 44.16, 44.14, 44.15, and 44.12, respectively, which are quite acceptable. The parameters used are $B=160$, $R=10$, and $s=(160/10)/2$. Therefore, the sizes of the authentication and recovery blocks are 160×160 and 10×10 , respectively. The distance between a recovery block and its corresponding embedding one is at least half of the height of the authentication block. Figure 8(c) and 8(c) show the tampered images and the results of tamper detection. The modifications are as follows. Airplane: the number 16 is changed to 10; Baboon: a polygon is place on its nose; Lena: a dot is placed on her face; and Peppers: a square is placed on one of the peppers. As can be seen, all of the tampered pixels are correctly identified, together with only a few false positives, which are shown in Figure 8(e). The number on the airplane is correctly detected with a few false positives. The dot on *Lena's* face is correctly detected without any false positive. The polygon is correctly detected and a convex shape is formed together with the false positives. Such result is identical to the previous analysis. And the square on the pepper is correctly detected with a few false positives. The recovered results are shown in Figure 8(f) and their PSNR are 43.64, 43.94, 44.15, and 44.11, respectively. We can see the tampered regions are well recovered.



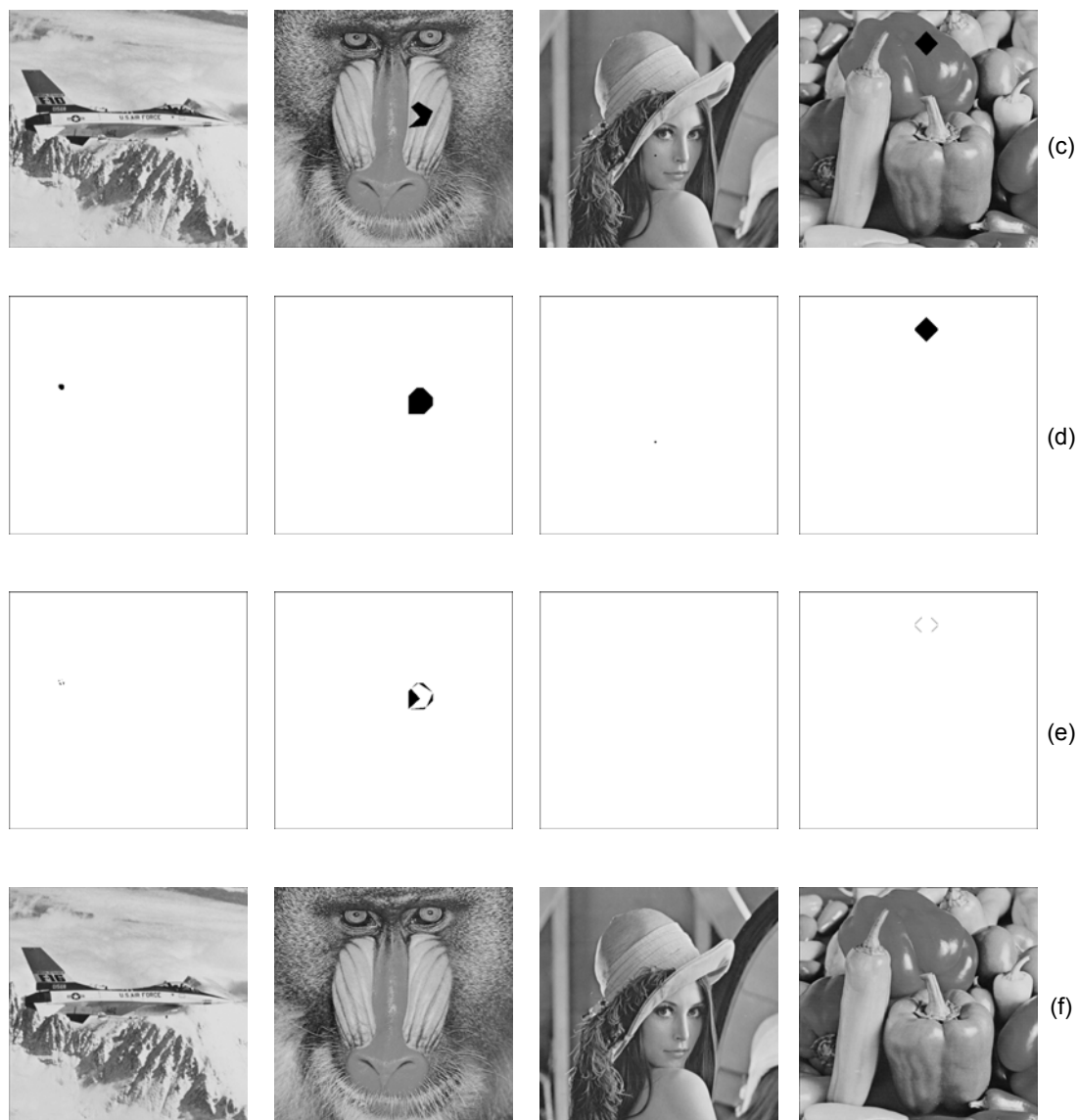


Figure 8: (a) original images, (b) data embedded images, (c) tampered images, (d) results of tamper detection, (e) false positives, and (f) recovered results of tampered images

4. Security Analysis

As the DES encryption system is used to generate the signatures, the proposed method is secure against the attack of manipulating individual image pixels. Three other attacks are the *search*, *collage*, and *cut-and-paste* attacks [9],[10], which are common for block-wise content authentication techniques. Because the attacked image has to maintain good visual quality, the size of the pasted blocks has to be very small in order for keeping the homogeneity of the block content. Therefore, the key requirement for these kinds of attacks to be successful is that the block size is small enough, usually less than or equal to 8×8 pixels. In the proposed method, as the size of the block is 160×160 , it clearly makes these attacks infeasible. That is, even if the attacker may forge an authentic image from a database containing hundreds of thousands of authentic images, it will certainly have poor visual quality due to block effects and incorrect block content. In conclusion, our method is invulnerable to these attacks.

5. Comparison of Detection Granularity

The detection granularity of the proposed method is compared with those of Patra *et al.*'s [4], Qi *et al.*'s [5], and Wu's [6] methods. Because the granularity of Qi *et al.*'s and Wu's methods depends on the image size, a unified size of 256×256 pixels is used in the analysis here. The comparison is shown in Table 1, in which it is obvious that our method outperforms the others.

Table 1. Comparison of the detection granularity

The proposed	Patra <i>et al.</i> 's	Qi <i>et al.</i> 's	Wu's
1×1=1 pixel	8×8=64 pixels	8×8=64 pixels	45 pixels

6. Conclusion

In this paper, we have described a technique to identify tampered pixels in an image. It is based upon dividing the image into authentication blocks and arranging linear signature blocks in such a way that they intersect at every pixel. As a consequence, each pixel is protected by four signatures and such an arrangement makes our technique capable of, in the best case, pinpointing a single altered pixel. This technique preserves the perceptual similarity of the original and the watermarked images, and it is also secure against various possible attacks. Although false positives are likely to be reported if altered pixels are spread randomly throughout the image, an attacker seems to have no reason to randomize the alterations. Therefore, our method is very useful for protecting the contents of the images at the granularity of one pixel. Moreover, the tampered region can be well recovered with the embedded recover data.

References

- [1] Liu H, Lin J, Huang J. *Image authentication using content based watermark*. Proc. IEEE Int. Sym. Circuits and Systems. 2005; 4: 4014-4017.
- [2] Rawat S, Raman B. A chaotic system based fragile watermarking scheme for image tamper detection. *Int. J. Electron. Commun.* 2011; 65: 840-847.
- [3] Xi'an Z. *A semi-fragile digital watermarking algorithm in wavelet transform domain based on Arnold transform*. Proc. IEEE Int. Conf. Signal Process. 2008: 2217-2220.
- [4] Patra JC, Phua JE, Rajan D. *DCT domain watermarking scheme using Chinese Remainder Theorem for image authentication*. Proc. IEEE Int. Conf. Multimedia and Expo. 2010: 111-116.
- [5] Qi X, Xin X, Chang R. *Image authentication and tamper detection using two complementary watermarks*. Proc. IEEE Int. Conf. Image Process. 2009: 4257-4260.
- [6] Wu Y. *Tamper-Localization Watermarking with Systematic Error Correcting Code*. Proc. IEEE Int. Conf. Image Process. 2006: 1965-1968.
- [7] Chen CH, Tang YL, Hsieh WS, Hwang MS. Image Protection by Intersecting Signatures. *TELKOMNIKA Indonesian Journal of Electrical Engineering*. 2014; 12: 6386-6392.
- [8] Wallace GK. The JPEG Still Picture Compression Standard. *IEEE Trans. Consumer Electronics*. 1992; 38: xviii-xxxiv.
- [9] Holliman M, Memon N. Counterfeiting Attacks on Oblivious Block-wise Independent Invisible Watermarking Schemes. *IEEE Trans. Image Process.* 2000; 9: 432-441.
- [10] Barreto PSLM, Kim HY, Rijmen V. Toward secure public-key blockwise fragile authentication watermarking. *IEE Proc. Vision, Image and Signal Processing*. 2002; 149: 57-62.
- [11] Liu Q. An Adaptive Blind Watermarking Algorithm for Color Image. *TELKOMNIKA Indonesian Journal of Electrical Engineering*. 2013; 11: 302-309.
- [12] Tan X, Hu D. A Watermarking Method Based on Optimization Statistics. *TELKOMNIKA Indonesian Journal of Electrical Engineering*. 2013; 11: 4794-4802.
- [13] Neyman SN, Pradnyana INP, Sitohang B. A New Copyright Protection for Vector Map using FFT-based Watermarking. *TELKOMNIKA Telecommunication, Computing, Electronics and Control*. 2014; 12: 367-378.
- [14] MT Suryadi, Nurpeti E. Performance of Chaos-Based Encryption Algorithm for Digital Image. *TELKOMNIKA Telecommunication, Computing, Electronics and Control*. 2014; 12: 675-682.